

① Machine learning foundation - Flashlight week 1

Getting started with IPython Notebook Graphlab

SHIFT+Enter = execute the command and SFrame

Enter: creates a multi-line code (doesn't execute)

Cell + New Cell + Markdown: write code in markdown

• Python is a scripting language, easy to prototype things
L or ESC(ESCAPE) + M

None type variable: n = None

Advanced printing: print "Our float value is %.s.
Our int value is %s." %- (f, i) → variable names

else if: elif

counter += 1: counter = counter + 1

Lambda function (other way to define a function):

E.g. square = lambda x: x * x . square(3) = 9

SFrames: scalable data structures for dealing with large datasets. It goes directly in the disk, you don't need lots of memory. It is part of Graphlab Create package.

Dataframe: dataframe. head() = print the first few lines

Dataframe. tail() print the last few lines.

② Machine Learning Foundations - Washington

Week 1
week 2

GraphLab Canvas

- Canvas for data visualization. It can take any data structure.

`dataset.show()` = open a new web page that represents the data visually with some interesting info = histogram, mean, most common value ... You can play with the canvas

To show this canvas directly inside the notebook

`graphlab.canvas.set_target('ipython')`

↳ means python notebook

With that, you cannot see all but select a category inside the canvas. E.g. `dataset['column'].show(view="Graphical")` one subgraphic view in the canvas or `dataset['column'].show()`

or `dataset['column'] = dataset['column']`.

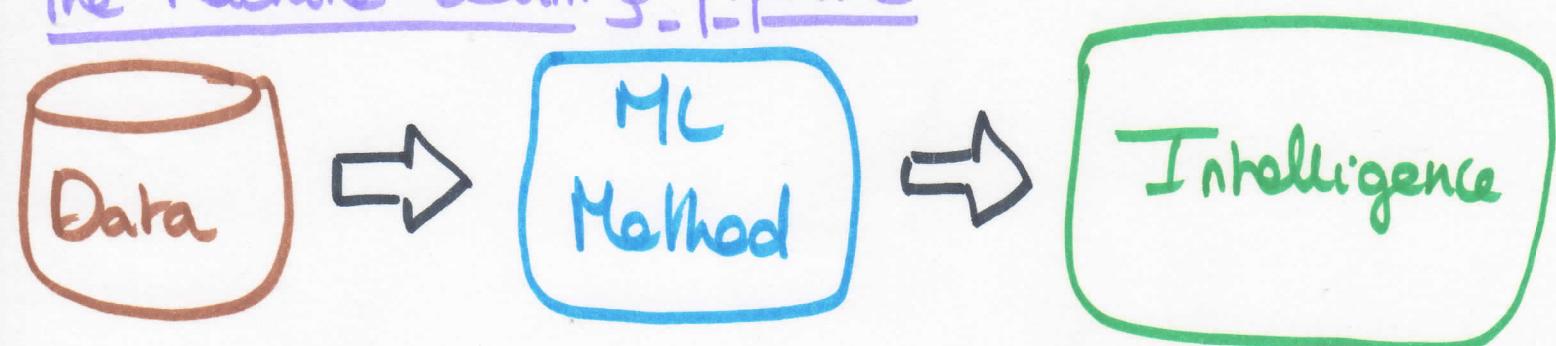
Apply function: `dataset['column'] = dataset['column'].apply(function)`

Regression

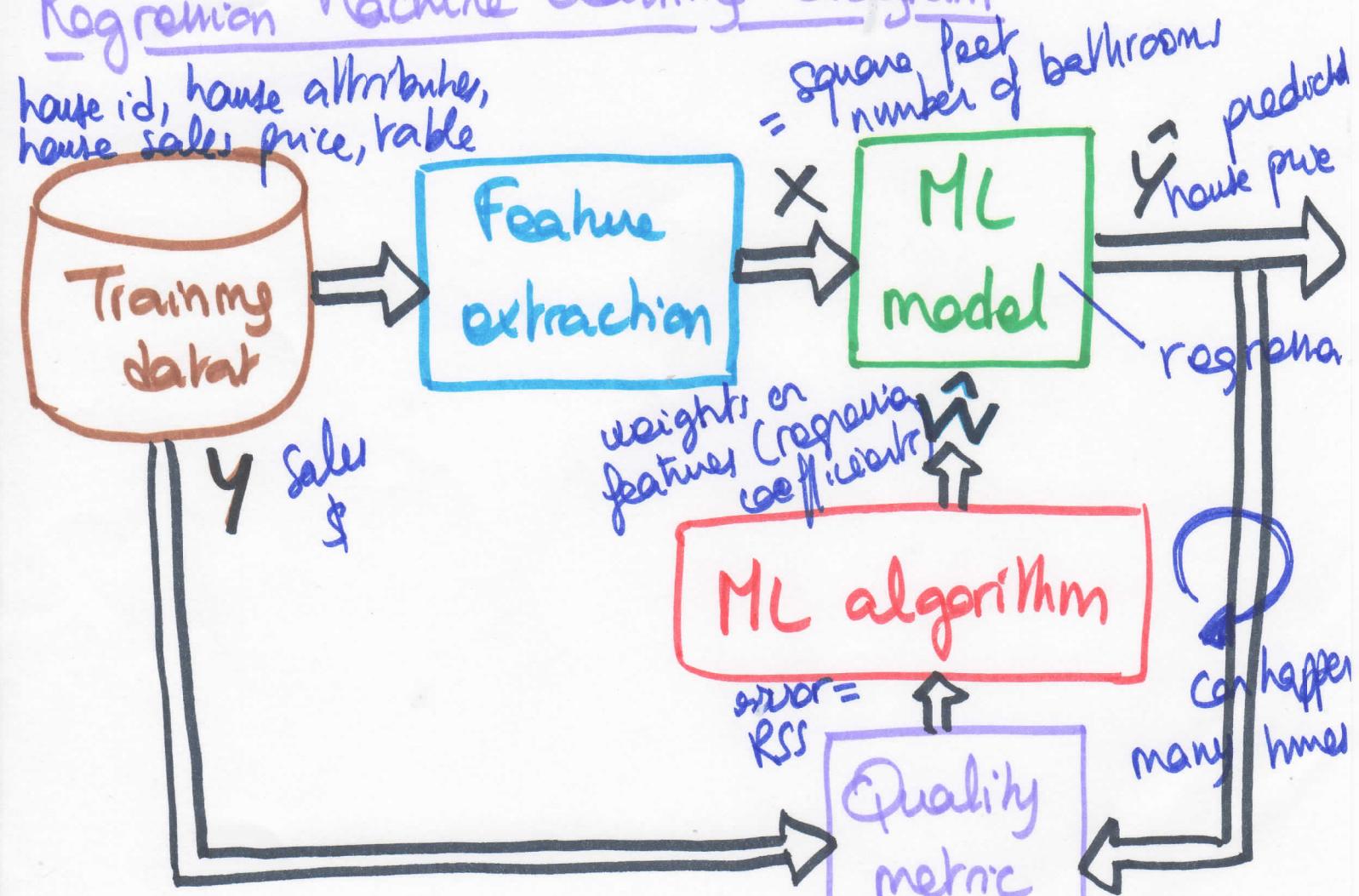
We have some set of features and we want to model how an observation that are associated with these features change as we change the values of the features. feature = covariate = predictor

③ Machine learning foundation - (Washington) week 1

The Machine Learning pipeline



Regression Machine Learning diagram



GraphLab house pricing (regression example)

To get some visualization: `dataset.show(view="Scatter Plot", x="column x", y="column y")` → click forget
`graphlab.canvas.set_target("ipynb")`

④ Machine learning - Foundations - (Washington) week 2

- You can get more details sliding the mouse on the data points and save the visualization as an image.

Splitting the data into training and test sets

```
train-data, test-data = dataset.random-split(  
.8, seed=2015)
```

↑ for 80% / 20% → random number

Build a simple regression model

```
model = graphlab.linear_regression.create(train-data,  
target = dependent variable, features = ['independent var'])
```

Evaluate the model (on the test data)

```
model.evaluate(test-data)  
{'max_error': ..., 'rmse': ...}
```

↑ outlier

Visualizing predictions with Matplotlib

```
import matplotlib.pyplot as plt  
%matplotlib inline
```

↑ piece of the package

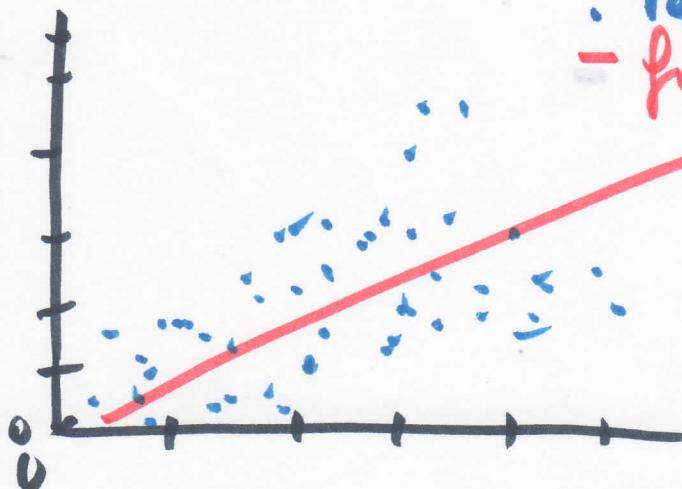
```
plot directly in the notebook  
plt.plot(test-data['independent var'], test-data['dependent var'], ':',  
         test-data['independent var'], model.predict(test-data), '-')
```

↑ scatterplot
↑ dot like

fitted
regression
line

⑤ Machine learning - foundations - Walking to week 2

- test data real value
- fitted regression line



Get the coefficients of the model

model.get('coefficients')

Exploring the features of the data

my-features = ["feature 1", "feature 2", ...]

dataset[my-features].show()

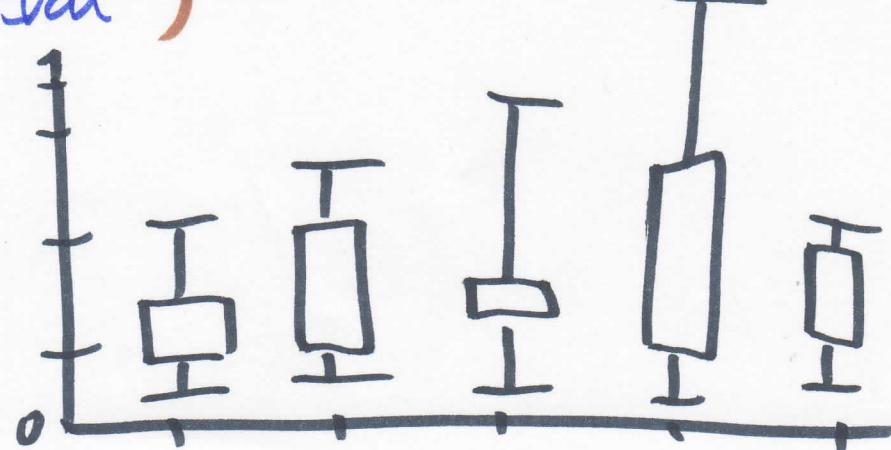
↳ you will get a nb of info on each feature (min, max, min, frequent items...)

Box plots

dataset.show(view="Box Whisker Plot", x="feature", y="dependent var")

↳ or independent var

• You can also play with the plot if lots of values



⑥ Machine Learning - Foundation - Washington week 2/3

Build a linear regression model with multiple features

my-features-model = graphlab.linear_regression.create(train-data, target='price', features='my-features')
↳ from previously

HTML in IPython notebook

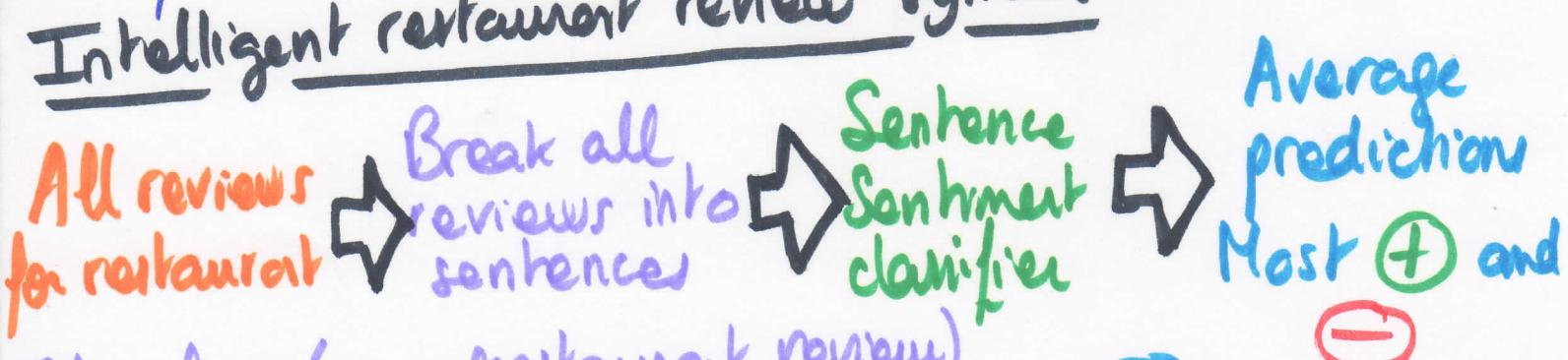
You can embed HTML, images ... in the notebook
 (e.g.) \approx markup

⚠ type ECHAP+M for the code to work

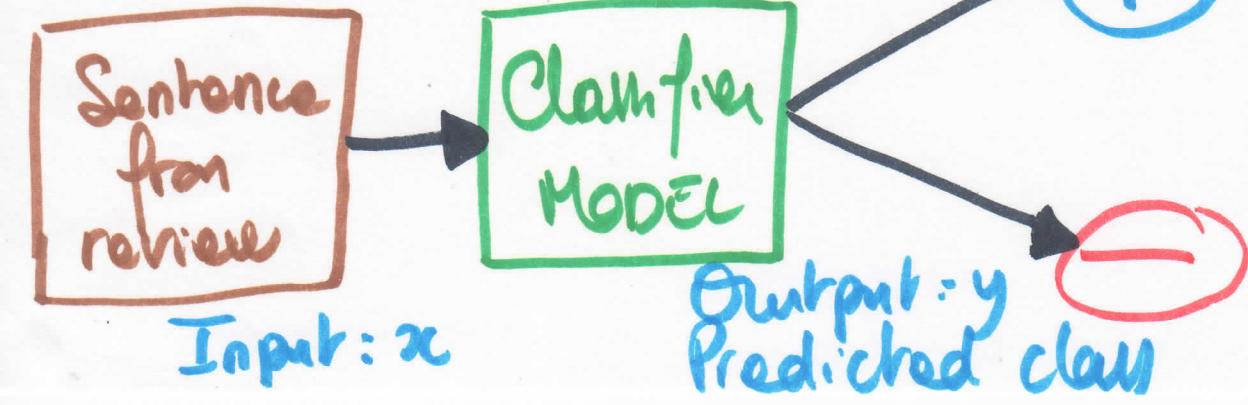
Classification

One of most common / used areas of Machine Learning.

Intelligent restaurant review system



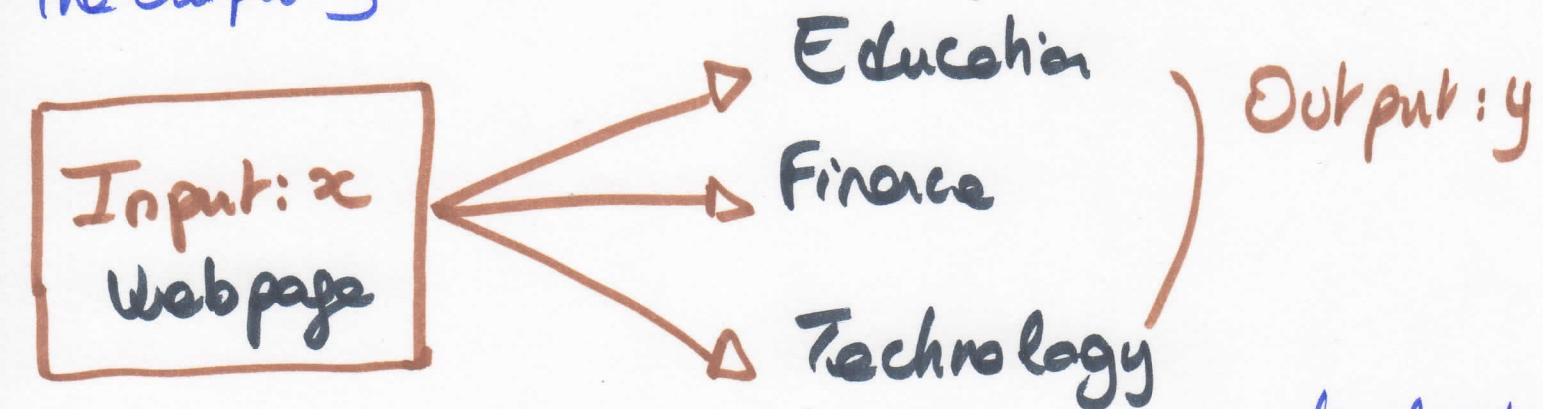
Classifier (e.g. restaurant review)



⑦ Machine learning foundations - Washington week 3

Multiclass classifier

The output y has more than 2 categories



- Classifier examples: spam filtering, image classification (through image pixels that recognize what it is), medical diagnosis (disease classifier)

Simple threshold classifier - restaurant review

List of positive words	List of negative words
great, awesome, good, amazing...	bad, terrible, disgusting,ucks...

Sushi was great, the food was awesome, but the service was terrible

Simple threshold classifier

Count positive & negative words in sentence

If $\frac{\text{number of positive words}}{\text{number of negative words}} > 1$ $\hat{y} = +$
Else $\hat{y} = -$

⑧ Machine learning - Foundations - Washington week 3

Problems with threshold classifier

- How do we get a list of positive/negative words?
- Words have different degrees of sentiment:
 - Great > good
 - How do we weight different words?
- Simple words are not enough
- Good \Rightarrow Positive
- Not good \Rightarrow Negative

Addressed by learning a classifier

Addressed by more elaborate features

Simple linear classifier

Word	Weight
good	1.0
great	1.2
:	:

• It is called a linear classifier, because the output is weighted sum of input.

Simple linear classifier

Score(x) = weighted count
of words in sentence

if Score(x) > 0 : $\hat{y} = +$

Else : $\hat{y} = -$

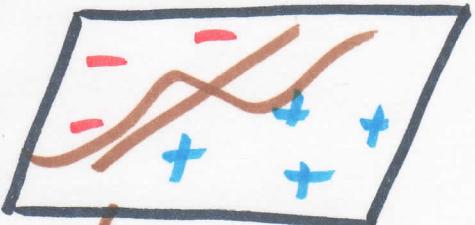
Sentence from review

Input: x

⑨ Machine learning - foundations - (Washington week 3)

Decision boundary separates positive & negative predictions

- For linear classifiers:
 - When 2 weights are non-zero: line
 - When 3 weights are non-zero: plane
 - When many weights are non-zero: hyperplane



- For more general classifiers
 - More complicated shapes

Classification error and accuracy (reviews example)

Error measures fraction of mistakes:

$$\text{error} = \frac{\# \text{ of mistakes}}{\text{Total } \# \text{ of sentences}}$$

Best possible value is 0.0

Measure accuracy - fraction of correct predictions

$$\text{accuracy} = \frac{\# \text{ of correct}}{\text{Total } \# \text{ of sentences}}$$

Best possible value is 1.0

What if you just guess?

$$\text{For binary classification: accuracy} = \underline{0.5}$$

$$\text{For k classes, accuracy} = \underline{1/k}$$

What is a good accuracy?

A classifier with 90% accuracy is not always good. E.g.
in spam email (to check ~)

⑩ Machine learning foundations - Washington week 3

- Always be digging in and asking the hard questions about reported accuracies.
 - Is there a class imbalance?
 - How does it compare to a simple baseline approach?
(random guessing, majority class...)
 - Most importantly = what accuracy does my application need? (What is good enough for my user's experience?
What is the impact of the mistakes we make?...)

Confusion matrix - False positives, false negatives

		Predicted label	
		+	-
True label	+	TP	FN
	-	FP	TN

- Cost of different types of mistakes can be different (α high) in some applications.

		Spam filtering	Medical diagnosis
False negative	Annoying	Disease not treated	Disease not treated → Useful treatment
	False positive	Email lost Higher cost	

- Multiclass classification

100 total examples

		Healthy	Cold	Flu
True label	Healthy	60	8	2
	Cold	12	4	4
True label	Flu	10	0	8
	Healthy	70	12	10

$$\begin{aligned} \text{accuracy} &= \\ &= 60 + 12 + 8 \\ &= 80 / 100 \end{aligned}$$

⑪ Machine learning foundations - Washington learning curve: How much data do I need? ^{week 3}

How much data does a model need to learn?

- **The more the merrier**

- But data quality is the most important factor 
Better less data but with more quality

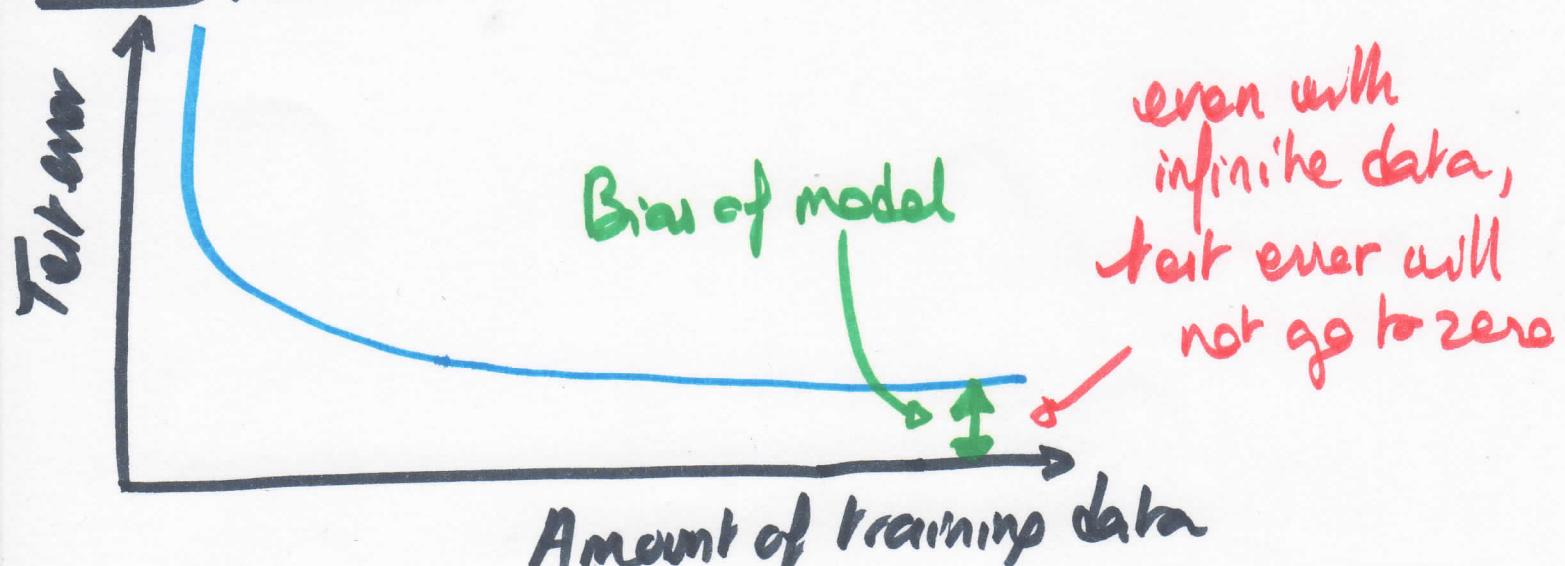
- **Theoretical techniques sometimes can bound how much data is needed**

- Typically too loose for practical application
- But provide guidance

- **In practice**

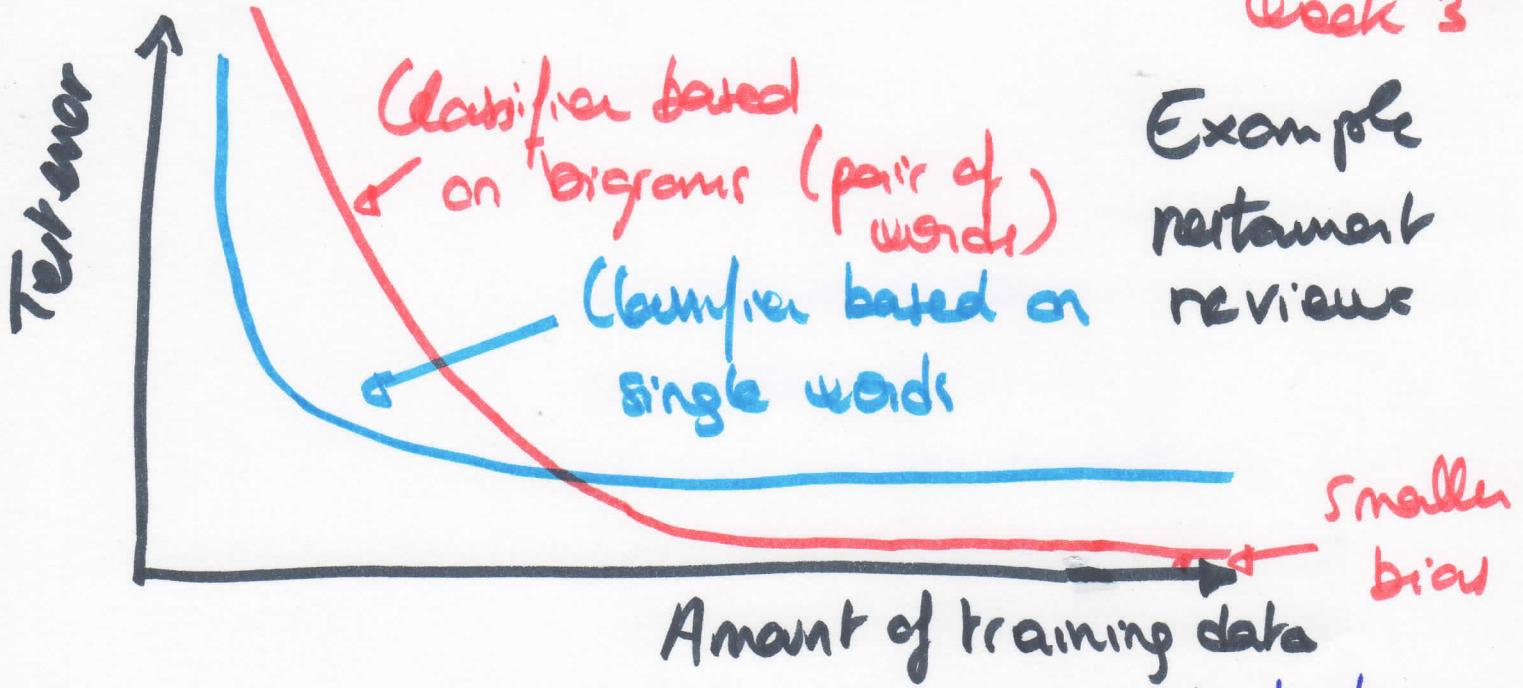
- More complex models require more data
- Empirical analysis can provide guidance

Learning curve



- More complex models tend to have less bias

⑫ Machine learning foundations - (Machine Learning Week 3)



- Models with less bias tend to need more data to learn well, but do better with sufficient data.

Class probabilities

- How confident is your prediction?
- Many classifiers provide a confidence level:

$$P(y=+ \text{ or } - | x) = ? \quad 0:1$$

Output label

Input sentence

+
or
-

E.g.:

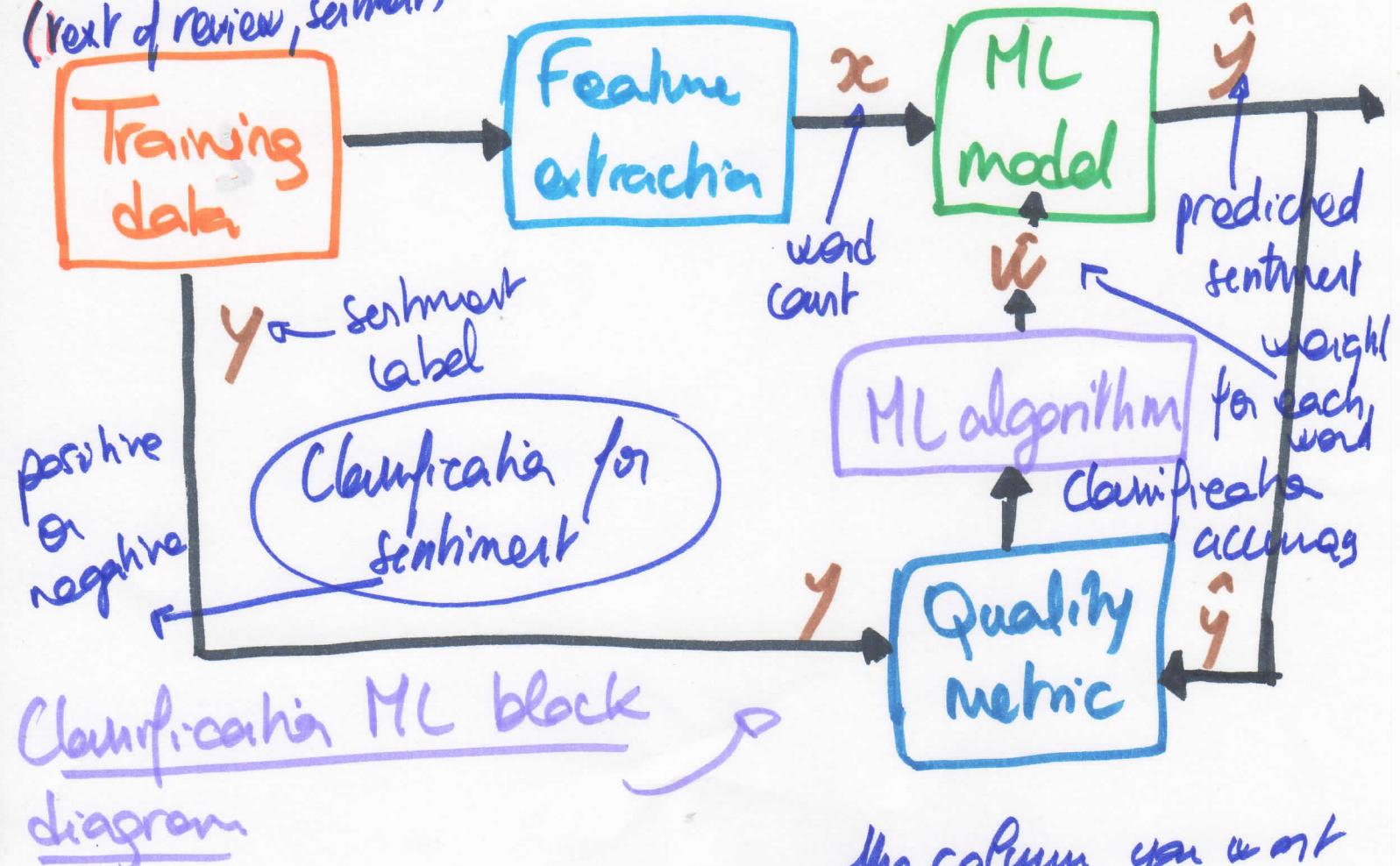
$$P(y=+ | x) = 0.99$$

High confidence

- It is extremely useful in practice.

13 Machine learning - foundations - (Washington, week 3)

(rest of review, sentiment)



Classification ML block diagram

Build a word count vector

the column you want to analyze.

`graphlab.text_analytics.count_words(dataset["column"])`

- It will return a list with pairs of word / number of times this word appears. (You can create a new column with this data - word count vector)

Build a classifier in graphlab python

`classifier_model = graphlab.logistic_classifier.create(train_data, target="column", features="column", validation_set= test_data)`

Add [] if the feature type is not a list (to convert it to a list)

14 Machine learning foundations - (Washington week 3/4)

Roc curve = evaluate of model

sentiment-model.evaluate(test-data, metric = "roc-curve")

- get some data, but for visualization with ability to change the threshold + get the confusion matrix:

sentiment-model.show("Evaluation")

Get the probabilities of the predictions

model.predict(dataset, output-type = "probability")

Get the last row / second last row... in a data frame

dataset[-1] dataset[-2] ...
get last row get second last row

Document retrieval

Currently reading an article you like.

Goal: you want to find similar article

Challenge: How do we measure similarity?

How do we search over articles?

Word count representation for measuring similarity

Bag of words model :- ignore order of words

- Count # of instances of each word in vocabulary

0 1 2 0 0 2 2 0 1 0 0 1

Carlos calls the sport
soccer. Tree calls sport cat football

"Carlos calls the sport
futbol. Emily calls the
sport soccer!"

(15) Machine Learning Foundations - Washington Week 4

Measuring similarity:

① Document talking about soccer

1)

0	0	0	0	5	3	0	0	1	0	0	0	0
2	0	0	0	2	0	0	1	0	1	1	c	c

$$1 \times 3 + 5 \times 2$$

$$= 13$$

similarities

② Document talking about soccer (some players)

2) ① Document talking about soccer

1	0	0	0	5	3	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---

$$= 0$$

no similarity

0	0	1	0	0	c	9	c	0	6	0	4	c
---	---	---	---	---	---	---	---	---	---	---	---	---

② Document talking about war in Africa

Issues with word counts: document length

The more the document is large the more the "similarity degree" will be high (as numbers will be higher for words)

1	0	0	c	1	5	3	1	0	0	1	0	c
---	---	---	---	---	---	---	---	---	---	---	---	---

Solution: normalize

$$\sqrt{1^2 + 5^2 + 3^2 + 1^2}$$

Prioritizing important words with tf-idf

Issues with word counts: rare words

Common words in doc: "the", "player", "field" .. dominate
 rare words like: "futbol", "Mehri"

16 Machine Learning Foundations - Washington week 4

Document frequency:

- What characterizes a rare word?
 - appears infrequently in the corpus
- Emphasize words appearing in fewer documents
 - equivalently, document word w based on # of docs containing w in corpus

Important words:

- Do we want only rare words to dominate?
- What characterizes an important word?
- Appear frequently in document (common locally)
 - Appear rarely in corpus (rare globally)
- Trade off between local frequency and global rarity

TF-IDF document representation

Term Frequency - Inverse Document Frequency (tf-idf)

Term frequency

1000		15	
The		Messi	

0			20
---	--	--	----

The
tf x idf Messi

Inverse document frequency 64 documents

10		4	
The		Messi	

$$\log\left(\frac{64}{1+63}\right) = 0 \quad \log\left(\frac{64}{1+63}\right) = \log(1) = 0$$

$$\log\left(\frac{\# \text{docs}}{1 + \# \text{docs using the word}}\right)$$

(17) Machine learning foundations - Washington
week 4

Idf log

- word in many docr $\log\left(\frac{\text{large \#}}{1+\text{large \#}}\right)$
- $= \log 1 = 0$
- rare word $\log\left(\frac{\text{large \#}}{1+\text{small \#}}\right) = \text{large \#}$

Retrieving similar documents

1- Nearest neighbor:

Input: Query article

Output: Most similar article (among the corpus)

Algorithm: Search over each article in corpus

- Compute $s = \text{similarity}(\text{query article}, \text{article in corpus})$
- if $s > Best-s$, record Most similar article = article in corpus
- and set $Best-s = s$
- Return Most similar article

K-Nearest neighbor

Input: Query article

Output: List of K similar articles

Nearest neighbor search: query article - corpus

specify: distance metric

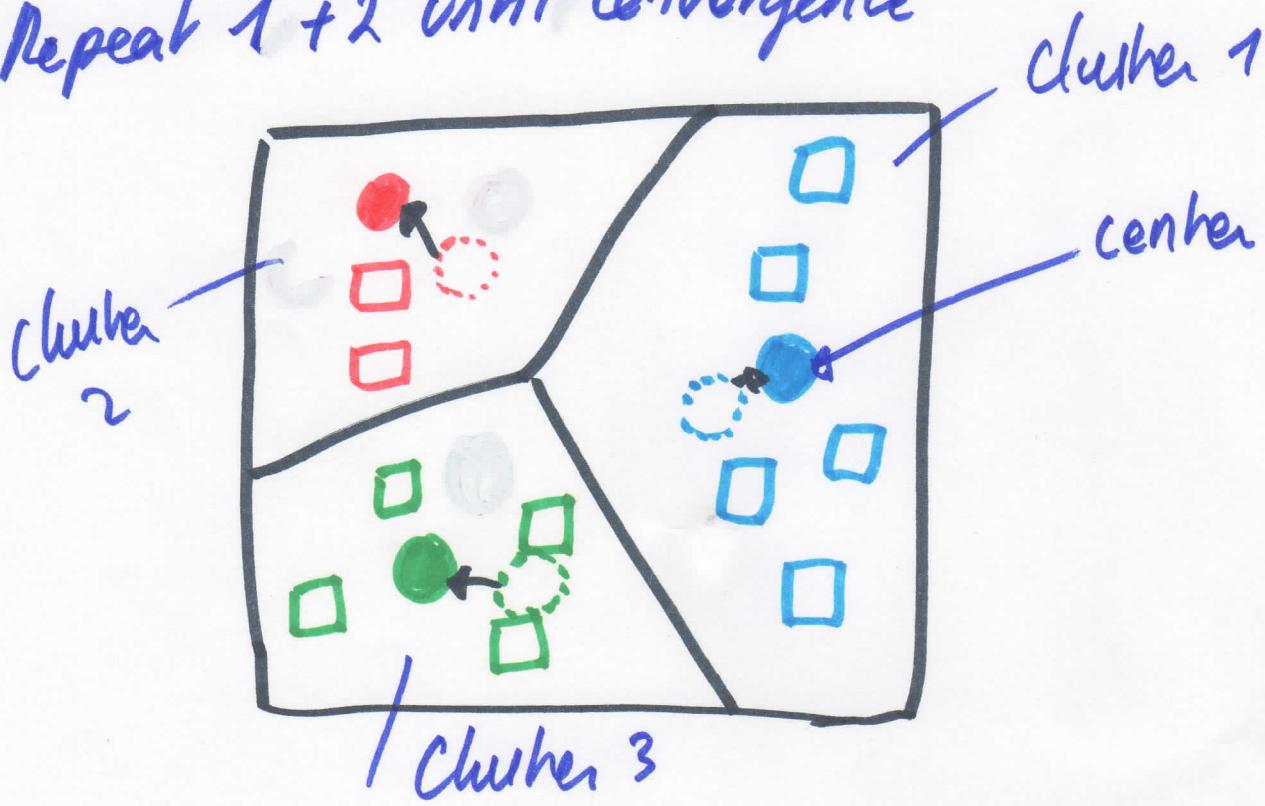
Output: Set of most similar articles

- (18) Machine learning foundations - Washington week 6
- Structure docs by topics: supervised learning
- Structure documents by topic. Discover groups (clusters) of related articles.
- Training set of labeled docs: sport, weather, science ..
- It is a multiclass classification problem
- Clustering documents: unsupervised learning Cluster 1
- No labels provided
 - Want to uncover cluster structure
- Input: docs as vectors
- Output: cluster labels
- Example with 2-dimensional vector
-
- word 1 word 2 word 3
- cluster 1 (label "sport" for instance)
- cluster 2
- cluster 3
- cluster 1 vector (doc)
- What defines a cluster?
- Cluster defined by center & shape / spread
 - Assign observation (doc) to cluster (topic label)
 - Score under cluster is higher than others
 - Often, just more similar to assigned cluster center than other cluster centers
- K-means algorithm
- Assume - similarity metric = distance to cluster center (smaller is better)

⑨ Machine learning foundations - Washington week 4

k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations
3. Repeat 1 + 2 until convergence

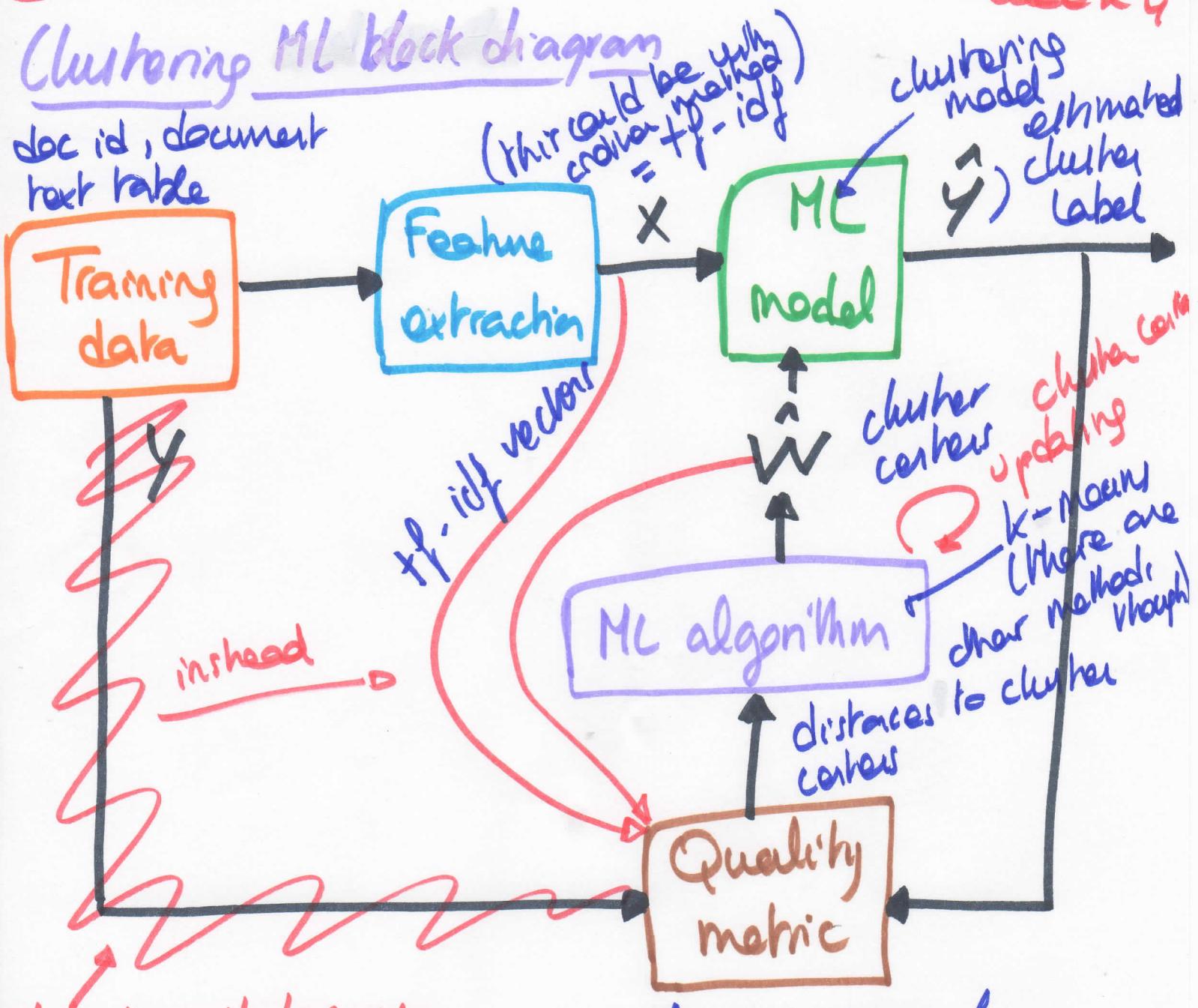


Other clustering examples

- Clustering images, for search, group as: ocean, dog...
- Grouping patients by medical condition (to better characterize subpopulation and disease)
- Products on Amazon (discover product categories from purchase histories, discovering groups of users)
- Structuring web search results (search terms can have multiple meanings).

⑩ Machine learning foundation - (Washington week 4)

Clustering ML block diagram



doesn't exist because unsupervised learning (we don't have any training data with labels)

Assess accuracy of our cluster labels

Goal: minimize distances between cluster centers and their associated cluster data points.

Graphlab SFrame print All rows and columns

SFrame.print_rows(num_rows=m, num_columns=n)

21) Machine learning - Foundation - Washington week 4

Python graphlab - Divide a column in two

Dictionary made of tuples \rightarrow column 1 = key, column 2 = value
dataset[['column name']].stack('column name', new-
column_name = ['column 1', 'column 2'])

Python graphlab = Compute TF-IDF

Have the column 'word-count' ready

tfidf = graphlab.text_analytics.tfidf(dataset[
'word-count'])

It returns a dictionary: You can then use .stack to divide
Tfidf into 2 columns: "word" and "tfidf", and sort
it by tfidf

Python graphlab = Compute distance TF-IDF

Manually, between two tfidf

graphlab.distance.cosine(dataset1['tfidf'][0],
dataset2['tfidf'][0]) to transform it into a dic-
tionary (instead of an array)

Cosine distance = the lower the cosine distance, the closer
the articles are. It measures similarity between vectors.
There are other techniques.

22) Machine Learning - Foundation - Washington weekly Python graphlab · build a k nearest neighbor model 5

knn-model = graphlab.nearest_neighbors.create(
dataset, features = ['tfidf'], label = 'name')
→ good for the example in the lecture, might be
tfidf column in the dataset with tfidf values

Query an input to find nearest

knn-model.query(input)

The input in the lecture video is the name of the person
E.g. arnold = people[people['name']] == 'Arnold Schwarzenegger'

Recommender systems = use earlier examples

Personalization is transforming our experience of the world

Youtube - video recommendation.

100 hours of videos are uploaded every minute.
What do I care about? There is information overload.
Browsing is "history": viewers | videos

Personalization: connects users & items

Movie recommendations - Netflix

Connect users with movies they may want to watch

Product recommendations - Amazon (the first to have built a recommender system)

23 Machine learning Foundations - Washington weeks

Recommendations combine global & user context

Music recommendations - Pandora

Recommendations from coherent & diverse sequence

Friend recommendations - Facebook

Users and "items" are of the same "type":

Build a network of users.

Drug-target interactions

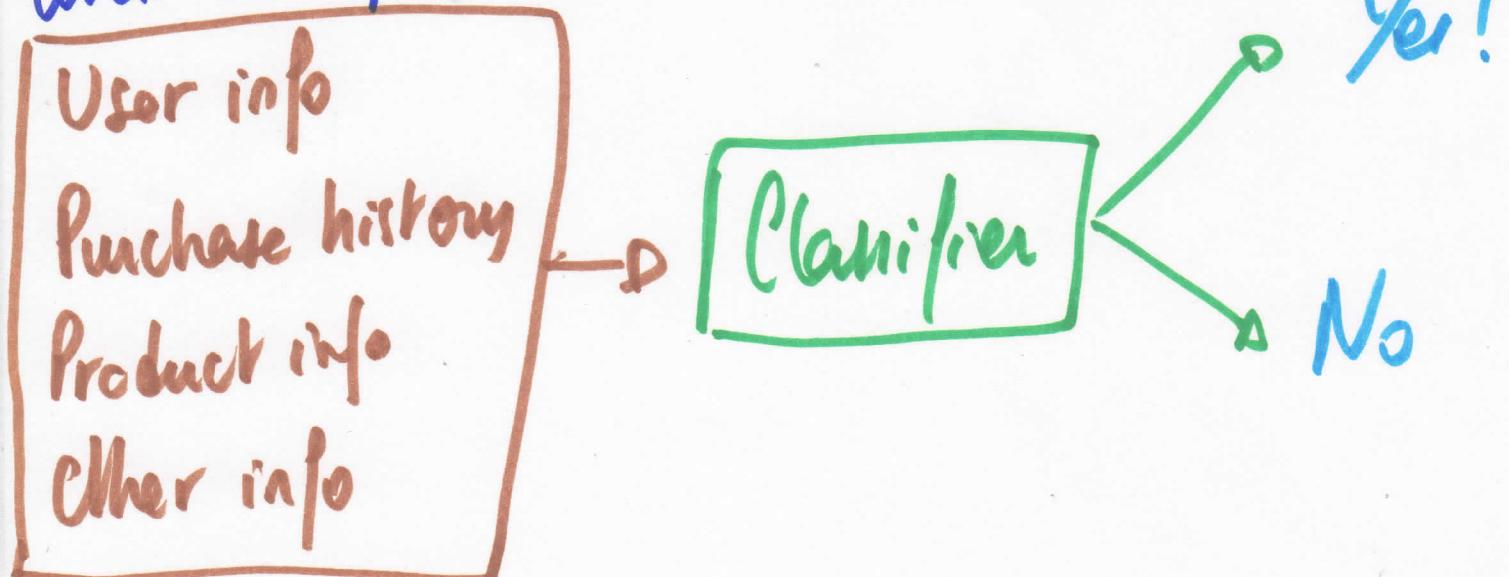
What drug should we "repurpose" for some disease?

Recommender system - Popularity (sol. 0)

- It is the simplest approach. Used by News websites
- What are people viewing now? Rank by global popularity.
- Limitations: no personalization "Most popular articles"

Recommender system - Classification model (sol. 1)

What is the probability I'll buy this product?



24 Machine learning foundations - Week 5

Pros:

- Personalized (consider user info & purchase history)
- Features can capture context: Time of the day
- What I just saw... Even handles limited user history: Ages of user...

Limitations:

- Features may not be available
- Often doesn't perform as well as collaborative filtering methods (next)

Recommender system = collaborative filtering (sol. 2)

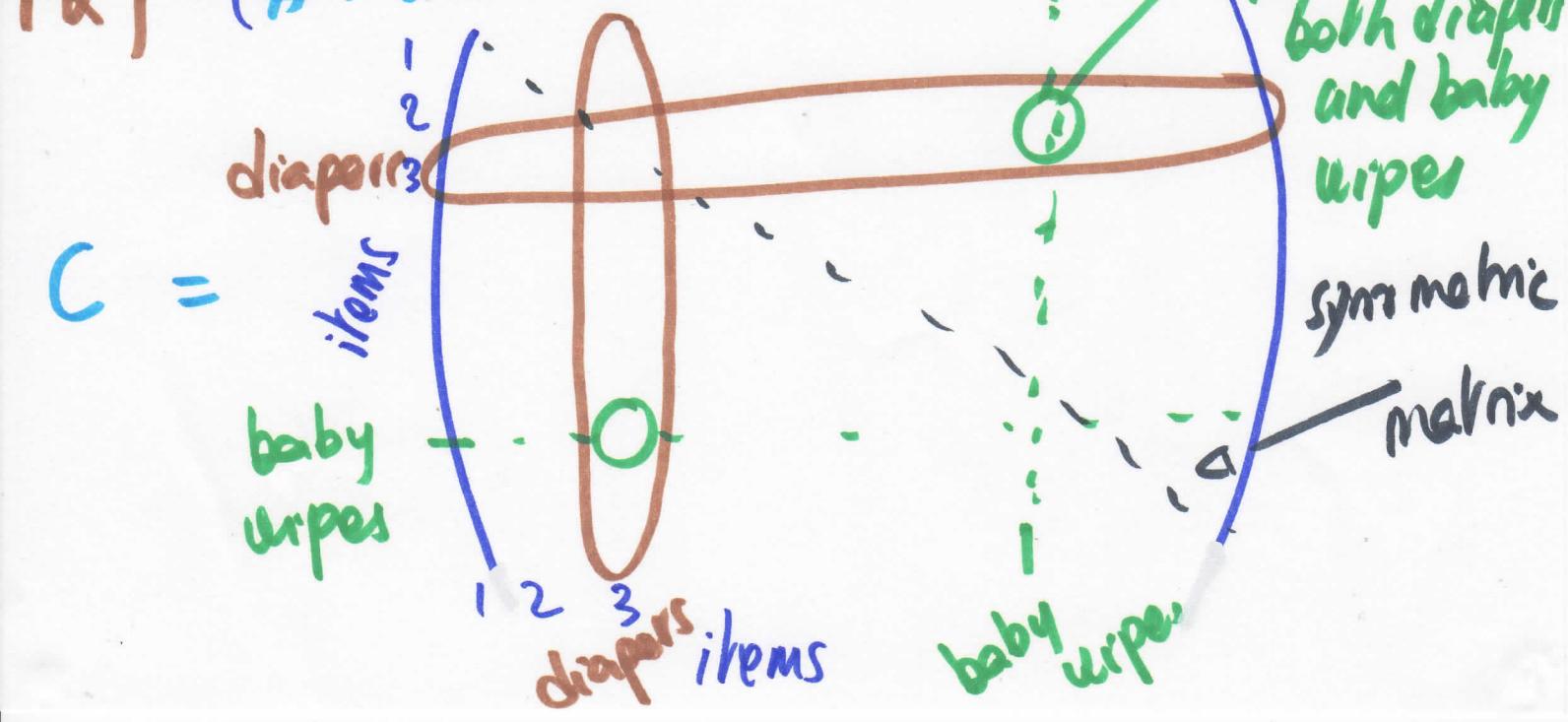
People who bought this, also bought....

Co-occurrence matrix - example

• People who bought diapers also bought baby wipes

• People who bought both items

Matrix C: store # users who bought both items
i & j - (# items x # items) matrix # of people purchasing both diapers and baby wipes



25) Machine learning foundation - Washington Week 5

Matrix C is symmetric: # purchasing i & j same
as # for $j \neq i$ ($C_{ij} = C_{ji}$)

Building this co-occurrence matrix: search for all
history users of purchase they've made and count. In
this example, every time we see a purchase of diaper
we'll add it to the matrix, and keep incrementing the
matrix this way as we are search through user history
of purchase.

Making recommendation using co-occurrences

User purchased diapers.

① look at diapers row of matrix
→ people who bought diapers
→ also bought all of
[... 0 ... 4 ... 100 ...] these products
DVD pacifier baby diapers

② Recommend other items with largest counts
(sorting the row) - here it could be baby diapers, milk.
for somebody who just purchased a product.

Co-occurrence matrix must be normalized

What if there are very popular items?

Popular baby item: diapers

Example: - For any baby item $i = \text{Sophie giraffe}$
large count C_{ij} for $j = \text{diaper}$

26 Machine learning foundation - Washington week 5

Sophie Giraffe [0 1million]
DVD diapers baby wipes

The recommender system will then be very likely to recommend diapers.

Results:

- Having these very large counts for popular items draws out all the other effects.
- We end up with a recommender system based on popularity \rightarrow need for normalization.

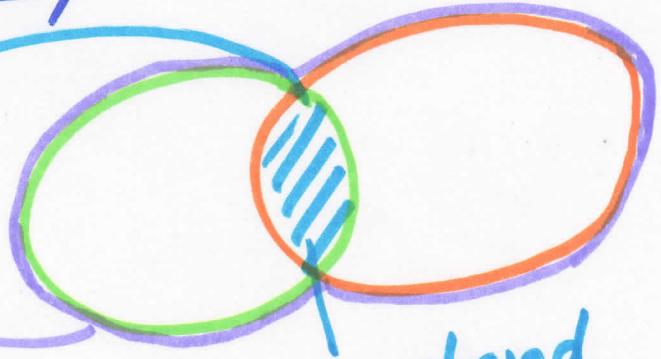
Normalize co-occurrences : similarity matrix

very similar to tf-idf

Jaccard similarity: normalizes by popularity
- Who purchased i and j divided by who purchased i or j

purchased i and j

purchased i or j



purchased items i and j

Many other similarity metrics possible, e.g. cosine similarity

Limitations

Only current page matter, no history

(27) Machine Learning Foundations - Washington week 5

- Recommend similar items to the one you bought
- What if you purchased many items?
 - We want a new approach that includes the purchase history

(Weighted) Average of purchased items

- User bought items {diapers, milk} → purchase history
- Compute user-specific score for each item j in inventory by combining similarities:

Score(user, baby wipes) =

$$1/2 (S_{\text{baby wipes, diapers}} + S_{\text{baby wipes, milk}})$$

some approach
but with milk rate

- Look at the rate diapers and see how many times people also bought baby wipes
- Average the results and decide how likely the user will buy baby wipes given his purchase history.
 - We could also weight recent purchases more and so on.

Sort Score(user, i) and find item j with highest similarity

Limitations: (well not directly)

- Does not utilize:
 - context (e.g. time of day) - product features
 - user features (e.g. age) (e.g. baby vs electronics)
 - because looking at all users

(28) Machine learning - foundation - Washington week 5

• Cold start problem

- What if a new user or product arrives? How do we form recommendation if we have absolutely no information on that product / user?

Discovering hidden structure by matrix factorization

(solution 3, recommender system)

Combines user history and product features and solves problems of new user / products.

Movie recommendation (context example)

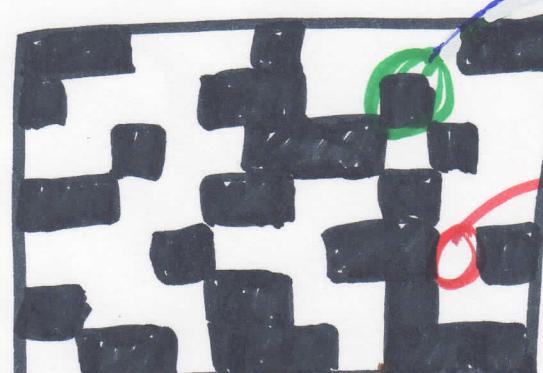
• Users watch movies and rate them

User	Movie	Rating
User 1	movie 1	*****
User 2	movie 2	****
User 3	movie 3	***
User 4	movie 4	*
User 5	movie 5	****
:	:	

- Each user only ^{has} watched a few of the available movies

Matrix completion problem (movie reco)

Rating =



User
v

Rating of movie v
by user u



Unknown
ratings

Movies

(29) Machine learning - Foundation - Washington week 5

Data: Users score same movies

Rating(u, v) known for black cells

Rating(u, v) unknown for white cells

Goal: Filling missing data? (white cells)

We use the known data to fill in missing data
(black cells) (white cells)

Predict the unknown ratings here

Suppose we had d features for each user and product

(movie example) (features)

Describe movie v with topics R_v

- How much is it action, romance, drama, ...

$$R_v = [0.3 \ 0.01 \ 1.5 \ \dots]$$

Describe user u with topics L_u

- How much she liked action, romance, drama, ...

$$L_u = [2.5 \ 0 \ 0.8 \ \dots]$$

The product of the two vectors

Rating(u, v) is the product of the two vectors

$$\text{estimate } R_v = [0.3 \ 0.01 \ 1.5 \ \dots] \rightarrow 0.3 \times 2.5 + 0 + 1.5 \times 0.8 + \dots = 7.2$$

$$L_u = [2.5 \ 0 \ 0.8 \ \dots] \rightarrow 0 + 0.01 \times 3.5 + 1.5 \times 0.01 \dots = 0.8$$

$$L_u' = [0 \ 3.5 \ 0.01 \ \dots]$$

③ Machine learning foundations - Washington week 5

Recomendation: sort movies user hasn't watched by Rating(u, v) estimate

Estimate
When the movie vector and the user vector agree a lot,
we will get a larger number than when they don't.
large rating

We will get a ~~higher rating~~
higher rating lower rating
We will sort all movies with their predicted rating and
recommend the ones with higher rating.
With this type of model we are not restricted to
keep the score (Rating estimator) between 0 and 5
(stars for movies), we just look the ones with higher scores
Result of

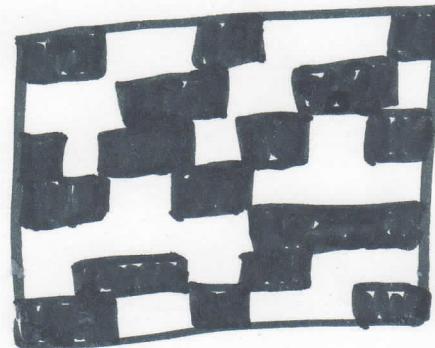
(stands for manner), we have
Predictions in matrix form

31 Machine learning foundations - (Washington week 5)

• But we don't know topics of users and movies...

Matrix factorization model: discovering ~~feature~~ from ~~data~~

We are going to estimate topics here
and L matrix (topics factors for every user and movie)
based on an observed rating (data)
 L and R represents the parameters (similar to regression
model that estimates parameters from data (black cells
in the rating matrix))



Rating =

